## Lecture 7
## Vector Processors &
## Multiprocessor Introduction

**Slides were used during lectures by
Krste Asanovic & David Patterson,
Berkeley, spring 2006**

---

## Outline

- Vector Processors
- Vector Metrics, Terms

- Multiprocessing Motivation
- SISD v. SIMD v. MIMD
- Centralized vs. Distributed Memory
- Challenges to Parallel Programming

- Conclusion

---

## Supercomputers

**Definition of a supercomputer:**
- Fastest machine in world at given task
- A device to turn a compute-bound problem into an I/O bound problem
- Any machine costing $30M+
- Any machine designed by Seymour Cray

**CDC6600 (Cray, 1964) regarded as first supercomputer**

---

## Supercomputer Applications

Typical application areas
- Military research (nuclear weapons, cryptography)
- Scientific research
- Weather forecasting
- Oil exploration
- Industrial design (car crash simulation)

**All involve huge computations on large data sets**

*In 70s-80s, Supercomputer ≡ Vector Machine*

---

## Vector Supercomputers

*Epitomized by Cray-1, 1976:*

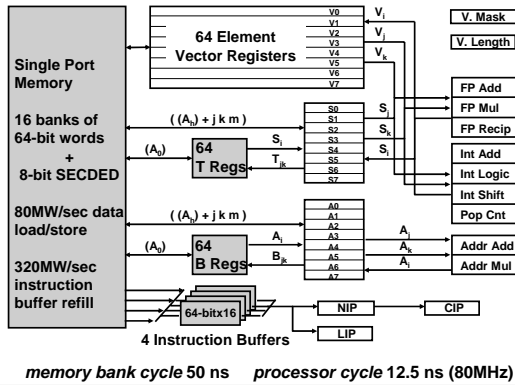**Scalar Unit + Vector Extensions**
- Load/Store Architecture
- Vector Registers
- Vector Instructions
- Hardwired Control
- Highly Pipelined Functional Units
- Interleaved Memory System
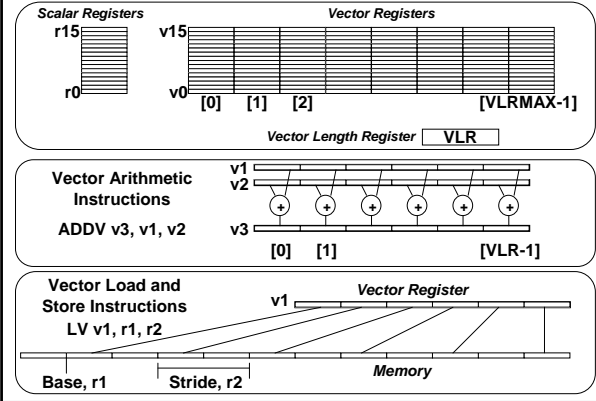- No Data Caches
- No Virtual Memory

---

## Cray-1 (1976)

## Cray-1 (1976)

**Single Port Memory**

**16 banks of 64-bit words + 8-bit SECDED**

**80MW/sec data load/store**

**320MW/sec instruction buffer refill**

64 Element Vector Registers

V0
V1
V2
V3
V4
V5
V6
V7

$V_i$
$V_j$
$V_k$

V. Mask

V. Length

$((A_h) + j k m)$

$S_i$

$S_j$
$S_k$
$S_i$

S0
S1
S2
S3
S4
S5
S6
S7

$(A_0)$ | 64 T Regs

$T_{jk}$

FP Add
FP Mul
FP Recip
Int Add
Int Logic
Int Shift
Pop Cnt

$((A_h) + j k m)$

$(A_0)$ | 64 B Regs

$A_i$
$B_{jk}$

A0
A1
A2
A3
A4
A5
A6
A7

$A_i$
$A_k$
$A_j$

Addr Add
Addr Mul

64-bitx16

NIP

CIP

LIP

**4 Instruction Buffers**

*memory bank cycle* **50 ns**    *processor cycle* **12.5 ns (80MHz)**

---

## Vector Programming Model

*Scalar Registers*
r15 ... r0

*Vector Registers*
v15 ... v0
[0] [1] [2] [VLRMAX-1]

*Vector Length Register* **VLR**

**Vector Arithmetic Instructions**

ADDV v3, v1, v2

v1
v2
(+) (+) (+) (+) (+) (+)
v3
[0] [1] [VLR-1]

**Vector Load and Store Instructions**

LV v1, r1, r2

*Vector Register*
v1

**Base, r1** | **Stride, r2** | *Memory*

---

## Vector Code Example

| # C code<br>for (i=0; i<64; i++)<br>  C[i] = A[i] + B[i]; | # Scalar Code<br>  LI R4, 64<br>loop:<br>  L.D F0, 0(R1)<br>  L.D F2, 0(R2)<br>  ADD.D F4, F2, F0<br>  S.D F4, 0(R3)<br>  DADDIU R1, 8<br>  DADDIU R2, 8<br>  DADDIU R3, 8<br>  DSUBIU R4, 1<br>  BNEZ R4, loop | # Vector Code<br>  LI VLR, 64<br>  LV V1, R1<br>  LV V2, R2<br>  ADDV.D V3, V1, V2<br>  SV V3, R3 |
|---|---|---|

---

## Vector Instruction Set Advantages

- **Compact**
  - one short instruction encodes N operations
- **Expressive, tells hardware that these N operations:**
  - are independent
  - use the same functional unit
  - access disjoint registers
  - access registers in the same pattern as previous instructions
  - access a contiguous block of memory (unit-stride load/store)
  - access memory in a known pattern (strided load/store)
- **Scalable**
  - can run same object code on more parallel pipelines or *lanes*

---

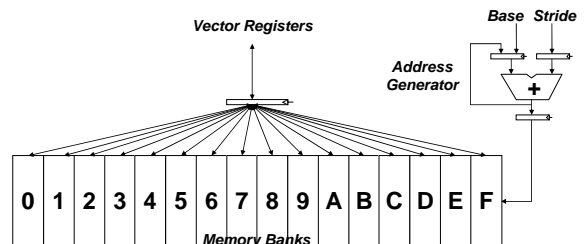## Vector Arithmetic Execution

- **Use deep pipeline ($\Rightarrow$ fast clock) to execute element operations**
- **Simplifies control of deep pipeline because elements in vector are independent ($\Rightarrow$ no hazards!)**
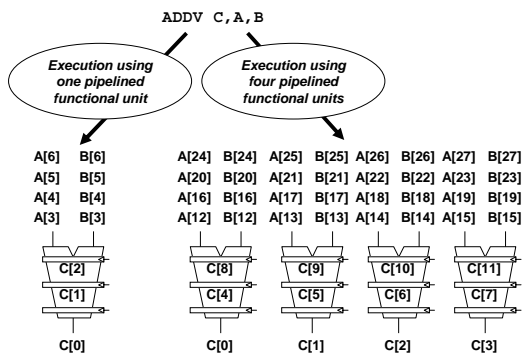
V1 V2 V3

*Six stage multiply pipeline*

**V3 <- V1 x V2**

---

## Vector Memory System

**Cray-1, 16 banks, 4 cycle bank busy time, 12 cycle latency**

- *Bank busy time*: Cycles between accesses to same bank
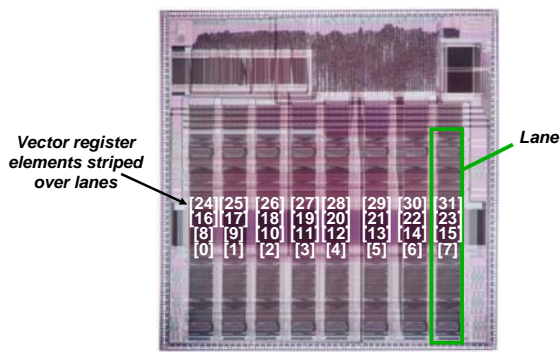
*Vector Registers*

**Base** **Stride**

*Address Generator*

(+)

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

*Memory Banks*

## Vector Instruction Execution

```
ADDV C,A,B
```

*Execution using one pipelined functional unit*

*Execution using four pipelined functional units*

| A[6] | B[6] |
|------|------|
| A[5] | B[5] |
| A[4] | B[4] |
| A[3] | B[3] |

C[2]
C[1]
C[0]

| A[24] | B[24] | A[25] | B[25] | A[26] | B[26] | A[27] | B[27] |
|-------|-------|-------|-------|-------|-------|-------|-------|
| A[20] | B[20] | A[21] | B[21] | A[22] | B[22] | A[23] | B[23] |
| A[16] | B[16] | A[17] | B[17] | A[18] | B[18] | A[19] | B[19] |
| A[12] | B[12] | A[13] | B[13] | A[14] | B[14] | A[15] | B[15] |

C[8]   C[9]   C[10]   C[11]
C[4]   C[5]   C[6]    C[7]
C[0]   C[1]   C[2]    C[3]

## Vector Unit Structure

*Functional Unit*

*Vector Registers*

| Elements 0, 4, 8, … | Elements 1, 5, 9, … | Elements 2, 6, 10, … | Elements 3, 7, 11, … |
|---|---|---|---|

*Lane*

*Memory Subsystem*

## T0 Vector Microprocessor (1995)

*Vector register elements striped over lanes*

*Lane*

[24] [25] [26] [27] [28] [29] [30] [31]
[16] [17] [18] [19] [20] [21] [22] [23]
[8]  [9]  [10] [11] [12] [13] [14] [15]
[0]  [1]  [2]  [3]  [4]  [5]  [6]  [7]

## Vector Memory-Memory versus Vector Register Machines

- **Vector memory-memory instructions hold all vector operands in main memory**
- **The first vector machines, CDC Star-100 ('73) and TI ASC ('71), were memory-memory machines**
- **Cray-1 ('76) was first vector register machine**

**Example Source Code**
```
for (i=0; i<N; i++)
{
  C[i] = A[i] + B[i];
  D[i] = A[i] - B[i];
}
```

**Vector Memory-Memory Code**
```
    ADDV C, A, B
    SUBV D, A, B
```

**Vector Register Code**
```
    LV V1, A
    LV V2, B
    ADDV V3, V1, V2
    SV V3, C
    SUBV V4, V1, V2
    SV V4, D
```
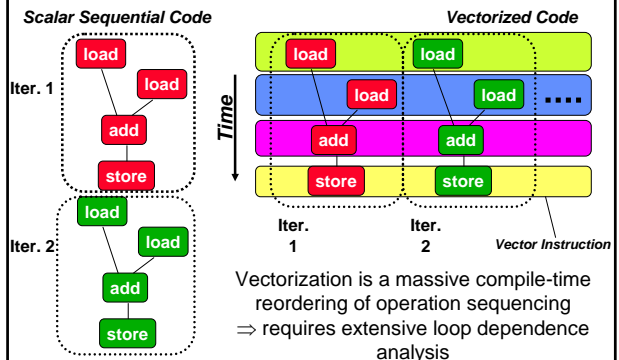
## Vector Memory-Memory vs. Vector Register Machines

- **Vector memory-memory architectures (VMMA) require greater main memory bandwidth, why?**
  - **All operands must be read in and out of memory**
- **VMMAs make if difficult to overlap execution of multiple vector operations, why?**
  - **Must check dependencies on memory addresses**
- **VMMAs incur greater startup latency**
  - Scalar code was faster on CDC Star-100 for vectors < 100 elements
  - For Cray-1, vector/scalar breakeven point was around 2 elements
- ⇒ *Apart from CDC follow-ons (Cyber-205, ETA-10) all major vector machines since Cray-1 have had vector register architectures*

  *(we ignore vector memory-memory from now on)*

## Automatic Code Vectorization

```
for (i=0; i < N; i++)
    C[i] = A[i] + B[i];
```

*Scalar Sequential Code*

*Vectorized Code*

*Time*

Iter. 1
load
load
add
store

Iter. 2
load
load
add
store

load
load
add
store

load
load
add
store

....

Iter. 1     Iter. 2

*Vector Instruction*

Vectorization is a massive compile-time reordering of operation sequencing ⇒ requires extensive loop dependence analysis

| Processor | Compiler | Completely vectorized | Partially vectorized | Not vectorized |
|---|---|---|---|---|
| CDC CYBER 205 | VAST-2 V2.21 | 62 | 5 | 33 |
| Convex C-series | FC5.0 | 69 | 5 | 26 |
| Cray X-MP | CFT77 V3.0 | 69 | 3 | 28 |
| Cray X-MP | CFT V1.15 | 50 | 1 | 49 |
| Cray-2 | CFT2 V3.1a | 27 | 1 | 72 |
| ETA-10 | FTN 77 V1.0 | 62 | 7 | 31 |
| Hitachi S810/820 | FORT77/HAP V20-2B | 67 | 4 | 29 |
| IBM 3090/VF | VS FORTRAN V2.4 | 52 | 4 | 44 |
| NEC SX/2 | FORTRAN77 / SX V.040 | 66 | 5 | 29 |

**Figure F.15 Result of applying vectorizing compilers to the 100 FORTRAN test kernels.** For each processor we indicate how many loops were completely vectorized, partially vectorized, and unvectorized. These loops were collected by Callahan, Dongarra, and Levine [1988]. Two different compilers for the Cray X-MP show the large dependence on compiler technology.

## Vector Stripmining

**Problem: Vector registers have finite length**

**Solution: Break loops into pieces that fit into vector registers, "Stripmining"**

```
for (i=0; i<N; i++)
  C[i] = A[i]+B[i];
```
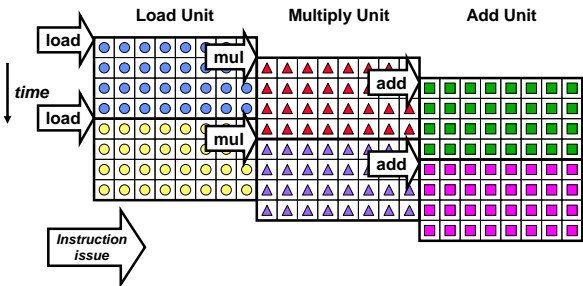


```
ANDI R1, N, 63    # N mod 64
MTC1 VLR, R1      # Do remainder
loop:
LV V1, RA
DSLL R2, R1, 3    # Multiply by 8
DADDU RA, RA, R2  # Bump pointer
LV V2, RB
DADDU RB, RB, R2
ADDV.D V3, V1, V2
SV V3, RC
DADDU RC, RC, R2
DSUBU N, N, R1    # Subtract elements
LI R1, 64
MTC1 VLR, R1      # Reset full length
BGTZ N, loop      # Any more to do?
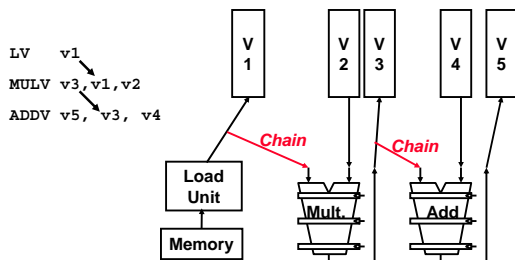```

## Vector Instruction Parallelism

**Can overlap execution of multiple vector instructions**
– example machine has 32 elements per vector register and 8 lanes



**Complete 24 operations/cycle while issuing 1 short instruction/cycle**

## Vector Chaining

- **Vector version of register bypassing**
  – introduced with Cray-1
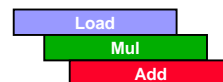
```
LV    v1
MULV v3,v1,v2
ADDV v5, v3, v4
```



## Vector Chaining Advantage

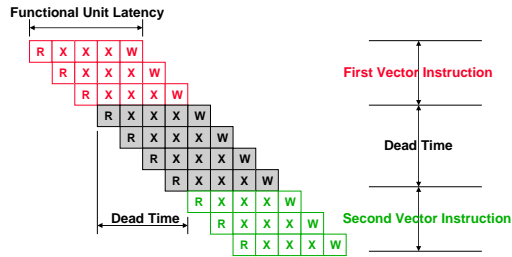- **Without chaining, must wait for last element of result to be written before starting dependent instruction**



- **With chaining, can start dependent instruction as soon as first result appears**
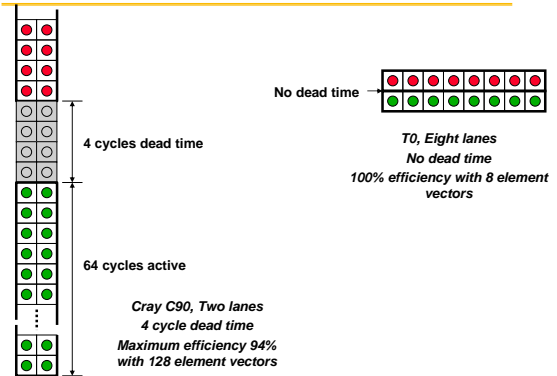
## Vector Startup

**Two components of vector startup penalty**
– functional unit latency (time through pipeline)
– dead time or recovery time (time before another vector instruction can start down pipeline)

Functional Unit Latency



First Vector Instruction

Dead Time

Dead Time

Second Vector Instruction

## Dead Time and Short Vectors



No dead time

4 cycles dead time

64 cycles active

*T0, Eight lanes*
*No dead time*
*100% efficiency with 8 element vectors*

*Cray C90, Two lanes*
*4 cycle dead time*
*Maximum efficiency 94% with 128 element vectors*

## Vector Scatter/Gather

**Want to vectorize loops with indirect accesses:**
```
for (i=0; i<N; i++)
    A[i] = B[i] + C[D[i]]
```

**Indexed load instruction (*Gather*)**
```
LV vD, rD        # Load indices in D vector
LVI vC, rC, vD   # Load indirect from rC base
LV vB, rB        # Load B vector
ADDV.D vA, vB, vC # Do add
SV vA, rA        # Store result
```

## Vector Scatter/Gather

**Scatter example:**
```
for (i=0; i<N; i++)
    A[B[i]]++;
```

**Is following a correct translation?**
```
LV vB, rB        # Load indices in B vector
LVI vA, rA, vB   # Gather initial A values
ADDV vA, vA, 1   # Increment
SVI vA, rA, vB   # Scatter incremented values
```

## Vector Conditional Execution

**Problem: Want to vectorize loops with conditional code:**
```
for (i=0; i<N; i++)
    if (A[i]>0) then
        A[i] = B[i];
```

**Solution: Add vector *mask* (or *flag*) registers**
– vector version of predicate registers, 1 bit per element
**…and *maskable* vector instructions**
– vector operation becomes NOP at elements where mask bit is clear
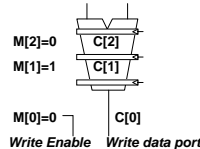
**Code example:**
```
CVM              # Turn on all elements
LV vA, rA        # Load entire A vector
SGTVS.D vA, F0   # Set bits in mask register where A>0
LV vA, rB        # Load B vector into A under mask
SV vA, rA        # Store A back to memory under mask
```
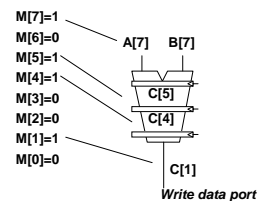
## Masked Vector Instructions

**Simple Implementation**
execute all N operations, turn off result writeback according to mask

**Density-Time Implementation**
scan mask vector and only execute elements with non-zero masks



M[7]=1  A[7]   B[7]
M[6]=0  A[6]   B[6]
M[5]=1  A[5]   B[5]
M[4]=1  A[4]   B[4]
M[3]=0  A[3]   B[3]

M[2]=0  C[2]
M[1]=1  C[1]

M[0]=0  C[0]
*Write Enable    Write data port*

M[7]=1
M[6]=0       A[7]   B[7]
M[5]=1
M[4]=1         C[5]
M[3]=0         C[4]
M[2]=0
M[1]=1        C[1]
M[0]=0
*Write data port*

## Compress/Expand Operations

- **Compress packs non-masked elements from one vector register contiguously at start of destination vector register**
  - population count of mask vector gives packed vector length
- **Expand performs inverse operation**



| M[7]=1 | A[7] | | A[7] | M[7]=1 |
| M[6]=0 | A[6] | | B[6] | M[6]=0 |
| M[5]=1 | A[5] | | A[5] | M[5]=1 |
| M[4]=1 | A[4] | | A[4] | M[4]=1 |
| M[3]=0 | A[3] | A[7] | B[3] | M[3]=0 |
| M[2]=0 | A[2] | A[5] | B[2] | M[2]=0 |
| M[1]=1 | A[1] | A[4] | A[1] | M[1]=1 |
| M[0]=0 | A[0] | A[1] | B[0] | M[0]=0 |

**Compress    Expand**

Used for density-time conditionals and also for general selection operations

## Vector Reductions

**Problem: Loop-carried dependence on reduction variables**
```
sum = 0;
for (i=0; i<N; i++)
    sum += A[i];  # Loop-carried dependence on sum
```

**Solution: Re-associate operations if possible, use binary tree to perform reduction**
```
# Rearrange as:
sum[0:VL-1] = 0                # Vector of VL partial sums
for(i=0; i<N; i+=VL)           # Stripmine VL-sized chunks
    sum[0:VL-1] += A[i:i+VL-1]; # Vector sum
# Now have VL partial sums in one vector register
do {
    VL = VL/2;                 # Halve vector length
    sum[0:VL-1] += sum[VL:2*VL-1] # Halve no. of partials
} while (VL>1)
```

## A Modern Vector Super: NEC SX-6 (2003)

- **CMOS Technology**
  - 500 MHz CPU, fits on single chip
  - SDRAM main memory (up to 64GB)
- **Scalar unit**
  - 4-way superscalar with out-of-order and speculative execution
  - 64KB I-cache and 64KB data cache
- **Vector unit**
  - 8 foreground VRegs + 64 background VRegs (256x64-bit elements/VReg)
  - 1 multiply unit, 1 divide unit, 1 add/shift unit, 1 logical unit, 1 mask unit
  - 8 lanes (8 GFLOPS peak, 16 FLOPS/cycle)
  - 1 load & store unit (32x8 byte accesses/cycle)
  - 32 GB/s memory bandwidth per processor
- **SMP structure**
  - 8 CPUs connected to memory through crossbar
  - 256 GB/s shared memory bandwidth (4096 interleaved banks)

## Multimedia Extensions

- **Very short vectors added to existing ISAs for micros**
- **Usually 64-bit registers split into 2x32b or 4x16b or 8x8b**
- **Newer designs have 128-bit registers (Altivec, SSE2)**
- **Limited instruction set:**
  - no vector length control
  - no strided load/store or scatter/gather
  - unit-stride loads must be aligned to 64/128-bit boundary
- **Limited vector register length:**
  - requires superscalar dispatch to keep multiply/add/load units busy
  - loop unrolling to hide latencies increases register pressure
- **Trend towards fuller vector support in microprocessors**

## Properties of Vector Processors

- **Each result independent of previous result**
  **=> long pipeline, compiler ensures no dependencies**
  **=> high clock rate**
- **Vector instructions access memory with known pattern**
  **=> highly interleaved memory**
  **=> amortize memory latency of over - 64 elements**
  **=> no (data) caches required! (Do use instruction cache)**
- **Reduces branches and branch problems in pipelines**
- **Single vector instruction implies lots of work (- loop)**
  **=> fewer instruction fetches**

## Operation & Instruction Count:
## RISC v. Vector Processor
(from F. Quintana, U. Barcelona.)

| Spec92fp | Operations (Millions) | | | Instructions (M) | | |
| Program | RISC | Vector | R / V | RISC | Vector | R / V |
| --- | --- | --- | --- | --- | --- | --- |
| swim256 | 115 | 95 | 1.1x | 115 | 0.8 | 142x |
| hydro2d | 58 | 40 | 1.4x | 58 | 0.8 | 71x |
| nasa7 | 69 | 41 | 1.7x | 69 | 2.2 | 31x |
| su2cor | 51 | 35 | 1.4x | 51 | 1.8 | 29x |
| tomcatv | 15 | 10 | 1.4x | 15 | 1.3 | 11x |
| wave5 | 27 | 25 | 1.1x | 27 | 7.2 | 4x |
| mdljdp2 | 32 | 52 | 0.6x | 32 | 15.8 | 2x |

Vector reduces ops by 1.2X, instructions by 20X

## Common Vector Metrics

- $R_\infty$: MFLOPS rate on an infinite-length vector
  - vector "speed of light"
  - Real problems do not have unlimited vector lengths, and the start-up penalties encountered in real problems will be larger
  - ($R_n$ is the MFLOPS rate for a vector of length n)
- $N_{1/2}$: The vector length needed to reach one-half of $R_\infty$
  - a good measure of the impact of start-up
- $N_V$: The vector length needed to make vector mode faster than scalar mode
  - measures both start-up and speed of scalars relative to vectors, quality of connection of scalar unit to vector unit

## Vector Execution Time

- Time = f(vector length, data dependicies, struct. hazards)
- *Initiation rate*: rate that FU consumes vector elements (= number of lanes; usually 1 or 2 on Cray T-90)
- *Convoy*: set of vector instructions that can begin execution in same clock (no struct. or data hazards)
- *Chime*: approx. time for a vector operation
- *m convoys take m chimes*; if each vector length is n, then they take approx. *m* x *n* clock cycles (ignores overhead; good approximization for long vectors)
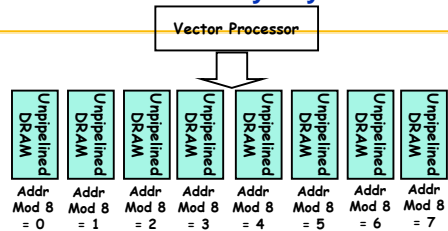
```
1: LV     V1,Rx     ;load vector X
2: MULV   V2,F0,V1  ;vector-scalar mult.
   LV     V3,Ry     ;load vector Y
3: ADDV   V4,V2,V3  ;add
4: SV     Ry,V4     ;store the result
```

4 convoys, 1 lane, VL=64
$\Rightarrow$ 4 x 64 = 256 clocks
(or 4 clocks per result)

## Memory operations

- Load/store operations move groups of data between registers and memory
- Three types of addressing
  - **Unit stride**
    » Contiguous block of information in memory
    » Fastest: always possible to optimize this
  - **Non-unit** (constant) **stride**
    » Harder to optimize memory system for all possible strides
    » Prime number of data banks makes it easier to support different strides at full bandwidth
  - **Indexed** (gather-scatter)
    » Vector equivalent of register indirect
    » Good for sparse arrays of data
    » Increases number of programs that vectorize

## Interleaved Memory Layout

Vector Processor

Unpipelined DRAM (×8)

Addr Mod 8 = 0 | Addr Mod 8 = 1 | Addr Mod 8 = 2 | Addr Mod 8 = 3 | Addr Mod 8 = 4 | Addr Mod 8 = 5 | Addr Mod 8 = 6 | Addr Mod 8 = 7

- **Great for unit stride:**
  - Contiguous elements in different DRAMs
  - Startup time for vector operation is latency of single read
- **What about non-unit stride?**
  - Above good for strides that are relatively prime to 8
  - Bad for: 2, 4
  - Better: prime number of banks…!

## How to get full bandwidth for Unit Stride?

- Memory system must sustain (# lanes x word) /clock
- No. memory banks > memory latency to avoid stalls
  - *m* banks $\Rightarrow$ *m* words per memory latency *l* clocks
  - if *m* < *l*, then gap in memory pipeline:

```
clock:  0 …  l   l+1  l+2 …   l+m-1   l+m  …  2l
word:   -- … 0    1    2  …   m-1      --  …   m
```

  - may have 1024 banks in SRAM
- If desired throughput greater than one word per cycle
  - Either more banks (start multiple requests simultaneously)
  - Or wider DRAMS. Only good for unit stride or large data types
- More banks/weird numbers of banks good to support more strides at full bandwidth
  - How to do prime number of banks efficiently?

## Vectors Are Inexpensive

### Scalar
- N ops per cycle $\Rightarrow$ O(N²) circuitry
- HP PA-8000
  - 4-way issue
  - reorder buffer: 850K transistors
    - incl. 6,720 5-bit register number comparators

### Vector
- N ops per cycle $\Rightarrow$ O(N + $\varepsilon$N²) circuitry
- T0 vector micro
  - 24 ops per cycle
  - 730K transistors total
    - only 23 5-bit register number comparators
  - No floating point

## Vectors Lower Power

### Single-issue Scalar
- One instruction fetch, decode, dispatch per operation
- Arbitrary register accesses, adds area and power
- Loop unrolling and software pipelining for high performance increases instruction cache footprint
- All data passes through cache; waste power if no temporal locality
- One TLB lookup per load or store
- Off-chip access in whole cache lines

### Vector
- One inst fetch, decode, dispatch per vector
- Structured register accesses
- Smaller code for high performance, less power in instruction cache misses
- Bypass cache
- One TLB lookup per group of loads or stores
- Move only necessary data across chip boundary

## Superscalar Energy Efficiency Even Worse

### Superscalar
- Control logic grows quadratically with issue width
- Control logic consumes energy regardless of available parallelism
- Speculation to increase visible parallelism wastes energy

### Vector
- Control logic grows linearly with issue width
- Vector unit switches off when not in use
- Vector instructions expose parallelism without speculation
- Software control of speculation when desired:
  - Whether to use vector mask or compress/expand for conditionals

## Vector Applications

*Limited to scientific computing?*
- **Multimedia Processing** (compress., graphics, audio synth, image proc.)
- **Standard benchmark kernels** (Matrix Multiply, FFT, Convolution, Sort)
- **Lossy Compression** (JPEG, MPEG video and audio)
- **Lossless Compression** (Zero removal, RLE, Differencing, LZW)
- **Cryptography** (RSA, DES/IDEA, SHA/MD5)
- **Speech and handwriting recognition**
- **Operating systems/Networking** (memcpy, memset, parity, checksum)
- **Databases** (hash/join, data mining, image/video serving)
- **Language run-time support** (stdlib, garbage collection)
- **even SPECint95**

## Older Vector Machines

| Machine | Year | Clock | Regs | Elements | FUs | LSUs |
|---|---|---|---|---|---|---|
| Cray 1 | 1976 | 80 MHz | 8 | 64 | 6 | 1 |
| Cray XMP | 1983 | 120 MHz | 8 | 64 | 8 | 2 L, 1 S |
| Cray YMP | 1988 | 166 MHz | 8 | 64 | 8 | 2 L, 1 S |
| Cray C-90 | 1991 | 240 MHz | 8 | 128 | 8 | 4 |
| Cray T-90 | 1996 | 455 MHz | 8 | 128 | 8 | 4 |
| Convex C-1 | 1984 | 10 MHz | 8 | 128 | 4 | 1 |
| Convex C-4 | 1994 | 133 MHz | 16 | 128 | 3 | 1 |
| Fuj. VP200 | 1982 | 133 MHz | 8-256 | 32-1024 | 3 | 2 |
| Fuj. VP300 | 1996 | 100 MHz | 8-256 | 32-1024 | 3 | 2 |
| NEC SX/2 | 1984 | 160 MHz | 8+8K | 256+var | 16 | 8 |
| NEC SX/3 | 1995 | 400 MHz | 8+8K | 256+var | 16 | 8 |

## Newer Vector Computers

- **Cray X1**
  - MIPS like ISA + Vector in CMOS
- **NEC Earth Simulator**
  - Fastest computer in world for 3 years; 40 TFLOPS
  - 640 CMOS vector nodes

## Key Architectural Features of X1

**New vector instruction set architecture (ISA)**
- Much larger register set (32x64 vector, 64+64 scalar)
- 64- and 32-bit memory and IEEE arithmetic
- Based on 25 years of experience compiling with Cray1 ISA

**Decoupled Execution**
- Scalar unit runs ahead of vector unit, doing addressing and control
- Hardware dynamically unrolls loops, and issues multiple loops concurrently
- Special sync operations keep pipeline full, even across barriers
- ⇒ Allows the processor to perform well on short nested loops

**Scalable, distributed shared memory (DSM) architecture**
- Memory hierarchy: caches, local memory, remote memory
- Low latency, load/store access to entire machine (tens of TBs)
- Processors support 1000's of outstanding refs with flexible addressing
- Very high bandwidth network
- Coherence protocol, addressing and synchronization optimized for DM
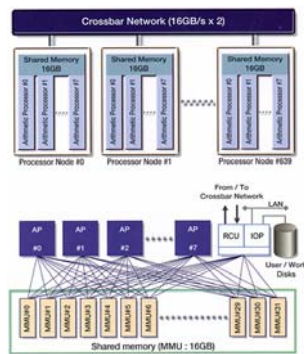
## Cray X1E Mid-life Enhancement

- **Technology refresh of the X1 (0.13μm)**
  - **~50% faster processors**
  - **Scalar performance enhancements**
  - **Doubling processor density**
  - **Modest increase in memory system bandwidth**
  - **Same interconnect and I/O**
- **Machine upgradeable**
  - **Can replace Cray X1 nodes with X1E nodes**
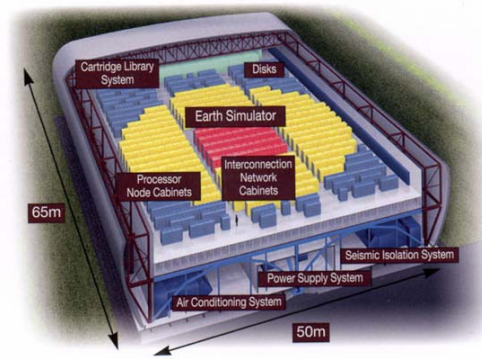
## ESS – configuration of a general purpose supercomputer

1. Processor Nodes (PN) Total number of processor nodes is 640. Each processor node consists of eight vector processors of 8 GFLOPS and 16GB shared memories. Therefore, total numbers of processors is 5,120 and total peak performance and main memory of the system are 40 TFLOPS and 10 TB, respectively. Two nodes are installed into one cabinet, which size is 40"x56"x80". 16 nodes are in a cluster. Power consumption per cabinet is approximately 20 KVA.

2. Interconnection Network (IN): Each node is coupled together with more than 83,000 copper cables via single-stage crossbar switches of 16GB/s x2 (Load + Store). The total length of the cables is approximately 1,800 miles.

3. Hard Disk. Raid disks are used for the system. The capacities are 450 TB for the systems operations and 250 TB for users.

4. Mass Storage system: 12 Automatic Cartridge Systems (STK PowderHorn9310); total storage capacity is approximately 1.6 PB.

*From Horst D. Simon, NERSC/LBNL, May 15, 2002, "ESS Rapid Response Meeting"*

## Earth Simulator



## Earth Simulator Building
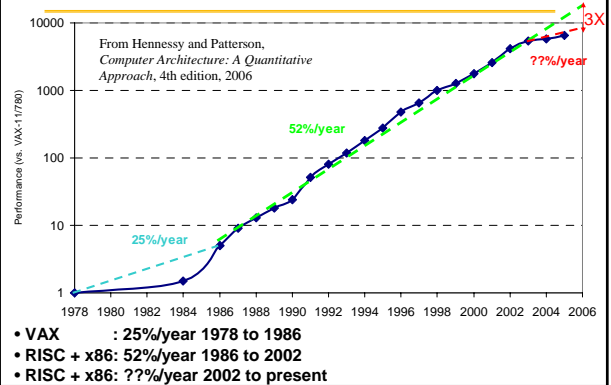


## ESS – complete system installed 4/1/2002



## Vector Summary

- **Vector is alternative model for exploiting ILP**
- **If code is vectorizable, then simpler hardware, more energy efficient, and better real-time model than Out-of-order machines**
- **Design issues include number of lanes, number of functional units, number of vector registers, length of vector registers, exception handling, conditional operations**
- **Fundamental design issue is memory bandwidth**
  - **With virtual address translation and caching**
- **Will multimedia popularity revive vector architectures?**

## Outline

- Vector Processors
- Vector Metrics, Terms

- Multiprocessing Motivation
- SISD v. SIMD v. MIMD
- Centralized vs. Distributed Memory
- Challenges to Parallel Programming

- Conclusion

## Uniprocessor Performance (SPECint)



From Hennessy and Patterson, *Computer Architecture: A Quantitative Approach*, 4th edition, 2006

- VAX       : 25%/year 1978 to 1986
- RISC + x86: 52%/year 1986 to 2002
- RISC + x86: ??%/year 2002 to present

## Déjà vu all over again?

**"… today's processors … are nearing an impasse as technologies approach the speed of light.."**
        David Mitchell, *The Transputer: The Time Is Now* (1989)
- Transputer had bad timing (Uniprocessor performance↑)
  ⇒ Procrastination rewarded: 2X seq. perf. / 1.5 years
- **"We are dedicating all of our future product development to multicore designs. … This is a sea change in computing"**
        Paul Otellini, President, Intel (2005)
- All microprocessor companies switch to MP (2X CPUs / 2 yrs)
  ⇒ Procrastination penalized: 2X sequential perf. / 5 yrs

| Manufacturer/Year | AMD/'05 | Intel/'06 | IBM/'04 | Sun/'05 |
|---|---|---|---|---|
| Processors/chip | 2 | 2 | 2 | 8 |
| Threads/Processor | 1 | 2 | 2 | 4 |
| Threads/chip | 2 | 4 | 4 | 32 |

## Other Factors ⇒ Multiprocessors

- **Growth in data-intensive applications**
  – Data bases, file servers, …
- **Growing interest in servers, server perf.**
- **Increasing desktop perf. less important**
  – Outside of graphics
- **Improved understanding in how to use multiprocessors effectively**
  – Especially server where significant natural TLP
- **Advantage of leveraging design investment by replication**
  – Rather than unique design

## Flynn's Taxonomy

M.J. Flynn, "Very High-Speed Computers", *Proc. of the IEEE*, V 54, 1900-1909, Dec. 1966.

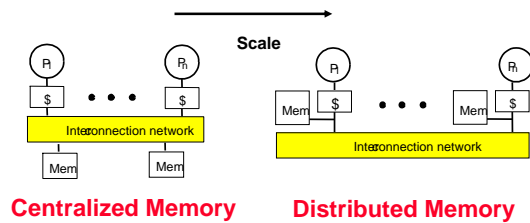- **Flynn classified by data and control streams in 1966**

| Single Instruction Single Data (SISD) (Uniprocessor) | Single Instruction Multiple Data SIMD (single PC: Vector, CM-2) |
|---|---|
| Multiple Instruction Single Data (MISD) (????) | Multiple Instruction Multiple Data MIMD (Clusters, SMP servers) |

- SIMD ⇒ Data Level Parallelism
- MIMD ⇒ Thread Level Parallelism
- MIMD popular because
  – Flexible: N pgms and 1 multithreaded pgm
  – Cost-effective: same MPU in desktop & MIMD

## Back to Basics

- **"A parallel computer is a collection of processing elements that underline{cooperate} and communicate to solve large problems fast."**
- **Parallel Architecture = Computer Architecture + Communication Architecture**
- **Two classes of multiprocessors WRT memory:**
  1. **Centralized Memory Multiprocessor**
     - < few dozen processor chips (and < 100 cores) in 2006
     - Small enough to share single, centralized memory
  2. **Physically Distributed-Memory multiprocessor**
     - Larger number chips and cores than 1
     - BW demands ⇒ Memory distributed among processors

## Centralized vs. Distributed Memory



**Centralized Memory**      **Distributed Memory**

## Centralized Memory Multiprocessor

- Also called <u>symmetric multiprocessors (SMPs)</u> because single main memory has a symmetric relationship to all processors
- Large caches $\Rightarrow$ single memory can satisfy memory demands of small number of processors
- Can scale to a few dozen processors by using a switch and by using many memory banks
- Although scaling beyond that is technically conceivable, it becomes less attractive as the number of processors sharing centralized memory increases

## Distributed Memory Multiprocessor

- Pro: Cost-effective way to scale memory bandwidth
  - If most accesses are to local memory
- Pro: Reduces latency of local memory accesses
- Con: Communicating data between processors more complex
- Con: Must change software to take advantage of increased memory BW

## Two Models for Communication and Memory Architecture

1. Communication occurs by explicitly passing messages among the processors: <u>message-passing multiprocessors</u>
2. Communication occurs through a shared address space (via loads and stores): <u>shared memory multiprocessors</u> either
   - **UMA** (Uniform Memory Access time) for shared address, centralized memory MP
   - **NUMA** (Non Uniform Memory Access time multiprocessor) for shared address, distributed memory MP
- In past, confusion whether "sharing" means sharing physical memory (Symmetric MP) or sharing address space

## Challenges of Parallel Processing

- First challenge is % of program inherently sequential
- Suppose 80X speedup from 100 processors. What fraction of original program can be sequential?
  - a. 10%
  - b. 5%
  - c. 1%
  - d. <1%

## Amdahl's Law Answers

$$\text{Speedup}_{\text{overall}} = \frac{1}{\left(1-\text{Fraction}_{\text{enhanced}}\right)+\dfrac{\text{Fraction}_{\text{parallel}}}{\text{Speedup}_{\text{parallel}}}}$$

$$80 = \frac{1}{\left(1-\text{Fraction}_{\text{parallel}}\right)+\dfrac{\text{Fraction}_{\text{parallel}}}{100}}$$

$$80\times\left(\left(1-\text{Fraction}_{\text{parallel}}\right)+\frac{\text{Fraction}_{\text{parallel}}}{100}\right)=1$$

$$79 = 80\times\text{Fraction}_{\text{parallel}}-0.8\times\text{Fraction}_{\text{parallel}}$$

$$\text{Fraction}_{\text{parallel}} = 79/79.2 = 99.75\%$$

## Challenges of Parallel Processing

- **Second challenge is long latency to remote memory**
- **Suppose 32 CPU MP, 2GHz, 200 ns remote memory, all local accesses hit memory hierarchy and base CPI is 0.5. (Remote access = 200/0.5 = 400 clock cycles.)**
- **What is performance impact if 0.2% instructions involve remote access?**
  - **a. 1.5X**
  - **b. 2.0X**
  - **c. 2.5X**

## CPI Equation

**CPI = Base CPI +**
  **Remote request rate x Remote request cost**

  **= 0.5 + 0.2% x 400 = 0.5 + 0.8 = 1.3**

**No communication is 1.3/0.5 or 2.6 faster than 0.2% instructions involve remote access**

## And in Conclusion [1/2] …

- **One instruction operates on vectors of data**
- **Vector loads get data from memory into big register files, operate, and then vector store**
- **E.g., Indexed load, store for sparse matrix**
- **Easy to add vector to commodity instruction set**
  - **E.g., Morph SIMD into vector**
- **Vector is very efficient architecture for vectorizable codes, including multimedia and many scientific codes**

## And in Conclusion [2/2] …

- **"End" of uniprocessors speedup => Multiprocessors**
- **Parallelism challenges: % parallalizable, long latency to remote memory**
- **Centralized vs. distributed memory**
  - **Small MP vs. lower latency, larger BW for Larger MP**
- **Message Passing vs. Shared Address**
  - **Uniform access time vs. Non-uniform access time**

## Reading and Schedule

- **This lecture:**
  - **Appendix F: *Vector Processors***
  - **Chapter 4: *4.1 Introduction Multiprocessors***

- **Next week, Oct 31st: *No class***

- **Next lecture, Nov 7th: *remainder of chapter 4* (in the afternoon *feedback* on assignment 2a)**

- **On Wed Nov 14th both at 11.15-13.00h and at 13.45-15.30h lectures in room 402**